

# **An Adaptive Partitioning Scheme for DRAM-based Cache in Solid State Drives**

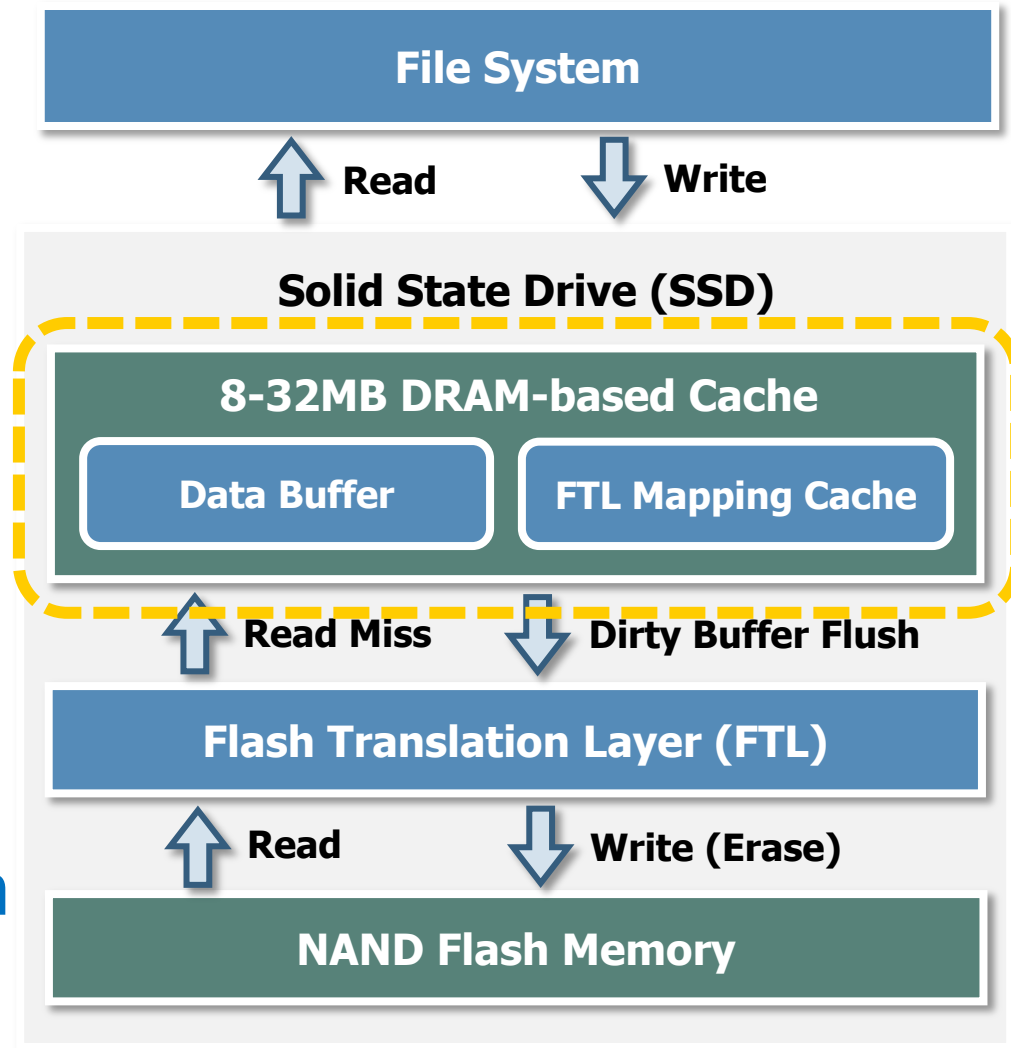
**Hyotaek Shim, Bon-Keun Seo, Jin-Soo Kim, and Seungryoul Maeng**

**Korea Advanced Institute of Science and Technology (KAIST)  
Sungkyunkwan University (SKKU)**

# Architecture of a Typical SSD

2/21

- The internal device cache has two main purposes
  - ▣ To absorb frequent read/write requests
  - ▣ To store logical-to-physical address mapping information
- We focus on **how to efficiently utilize the device cache between the two purposes**



# Trade-Off between Buffering and Mapping

3/21

- A trade-off between how much space is allocated to buffering versus mapping
  - ▣ Frequent requests can be more cached with larger buffering space
  - ▣ The SSD performance can also benefit from larger mapping space
- The device cache should be appropriately partitioned
  - ▣ Existing studies assumed that the *BM ratio* is fixed (*static partitioning policy*)

BM ratio is the ratio of the buffering and the mapping space.

8-32MB DRAM-based Cache

Data Buffer

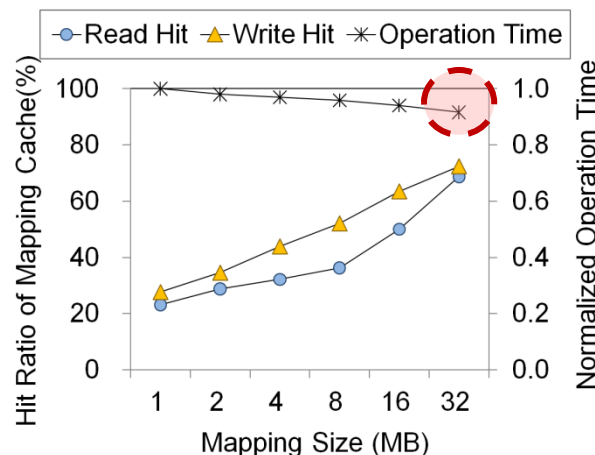
FTL Mapping Cache

# Effects of Adjusting the BM Ratio

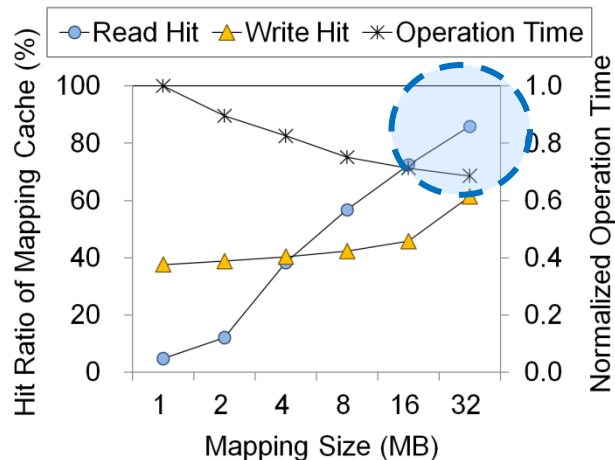
4/21

□ The optimal BM ratio is usually affected by workload characteristics

PC Trace



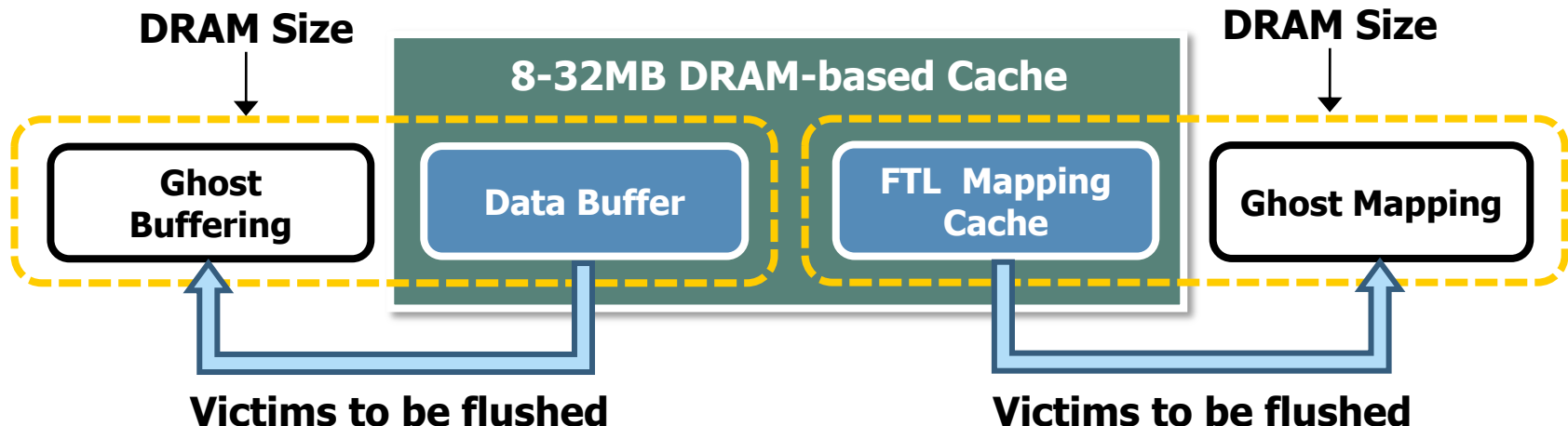
TPC-C Trace



# Adaptive Partitioning Scheme (1/3)

5/21

- This scheme adaptively adjusts the BM ratio according to workload characteristics
  - ▣ Comparing the **cost-benefits** of buffering and mapping
  - ▣ Ghost cache
    - Exclusive victim cache that stores only metadata
    - The **cost-benefits of actual caches are estimated by their ghost caches**



# Adaptive Partitioning Scheme (2/3)

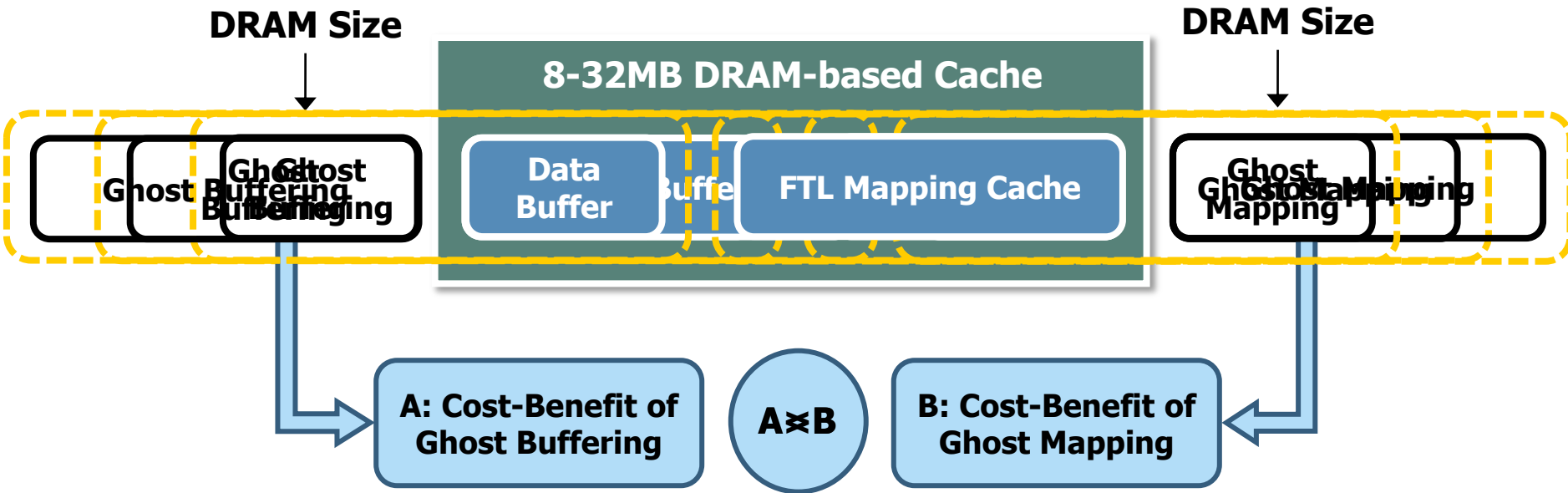
6/21

- ***Cost-benefit*** of ghost cache
  - ▣ Whenever a read/write hit occurs in ghost cache, its ***benefit*** is accumulated
    - At the same time, a read/write miss occurs in its actual cache
    - The ***benefit*** is the cost (**NAND flash operation time**) caused by the read/write miss in its actual cache
    - We call this cost **opportunity cost** caused by not enlarging the actual cache size
  - ▣ The ***cost*** of ***cost-benefit*** is the expected memory consumption of ghost cache

# Adaptive Partitioning Scheme (3/3)

7/21

- At every pre-defined interval, the BM ratio is tuned by comparing the cost-benefits of ghost caches

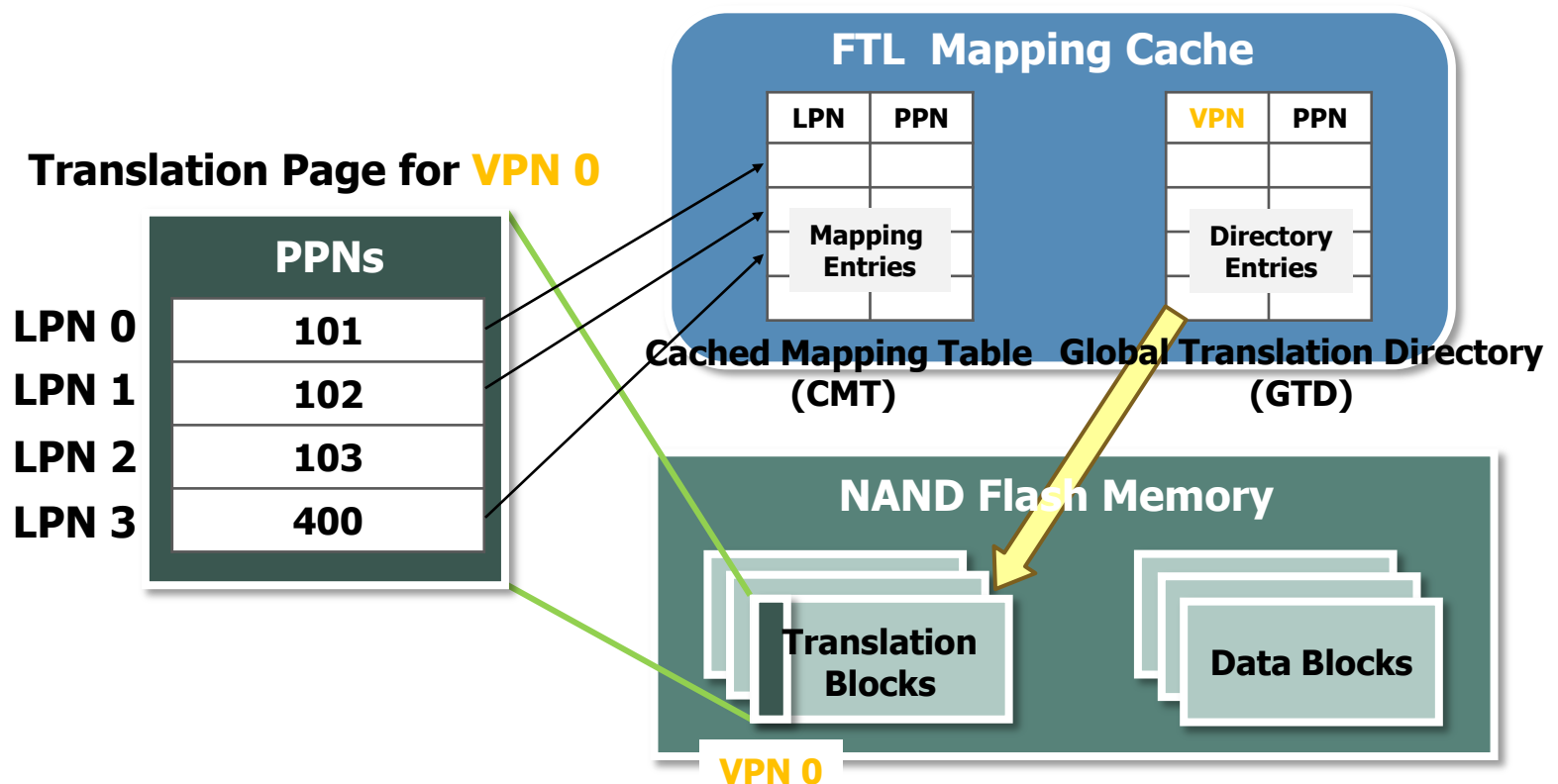


# Case Study I

## Demand-based Flash Translation Layer (DFTL)

8/21

- DFTL applies a caching mechanism to existing page-level mapping FTL
  - ▣ DFTL keeps only frequently-accessed logical-to-physical mapping entries in CMT

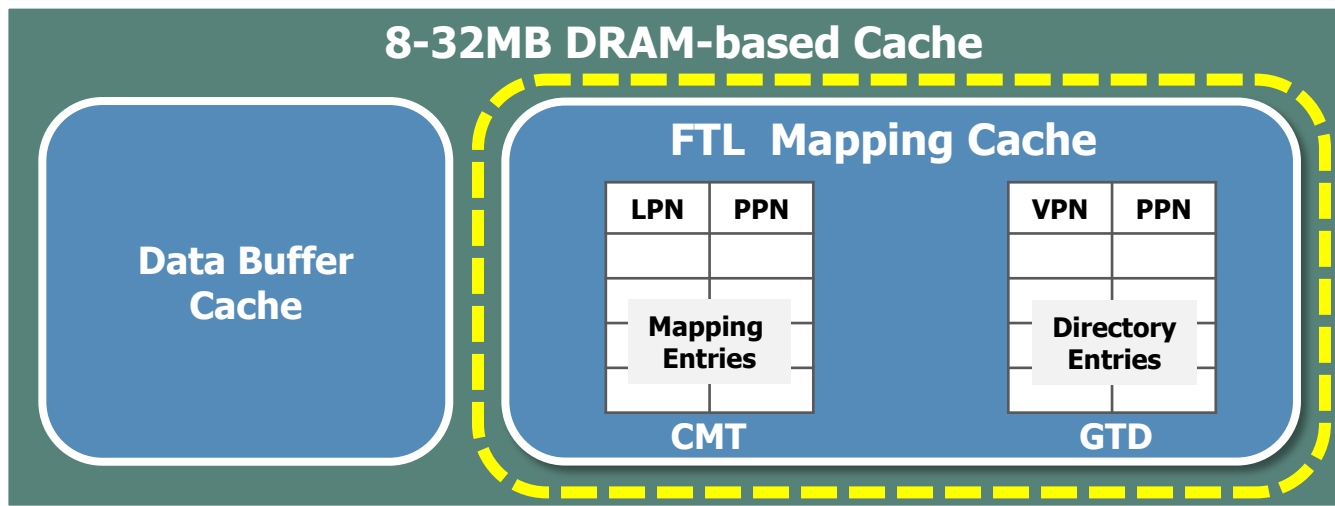




# Opportunity Cost of the Mapping Cache with DFTL

9/21

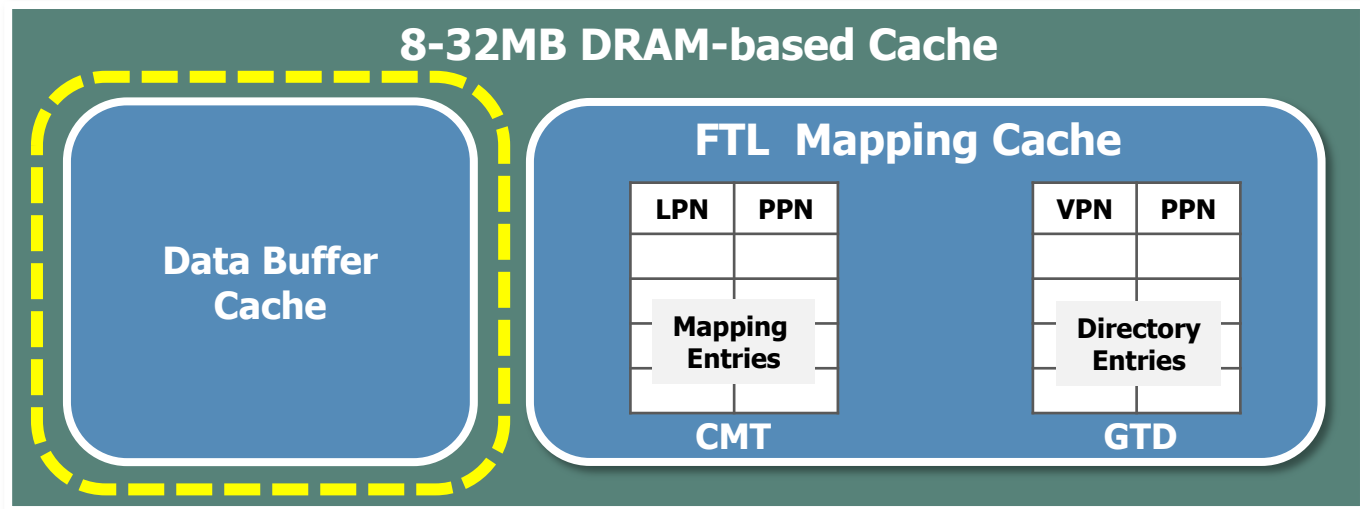
- Opportunity cost for a read miss in CMT
  - *Flash\_Read*
- Opportunity cost for a write miss in CMT
  - $\{Flash\_Read + Flash\_Write (GC\ Overhead)\} / Batch\_Factor$
  - *Batch\_Factor* means the avg. # of CMT entries flushed by a batch update



# Opportunity Cost of the **Buffer Cache** with DFTL

10/21

- **Opportunity cost for a read miss in buffer cache**
  - ▣ *Flash\_Read + Flash\_Read \* CMT\_Read\_Miss\_Ratio*
- **Opportunity cost for a write miss in buffer cache**
  - ▣ *Flash\_Write (GC\_Overhead)*  
+ *CMT\_Write \* CMT\_Write\_Miss\_Ratio*

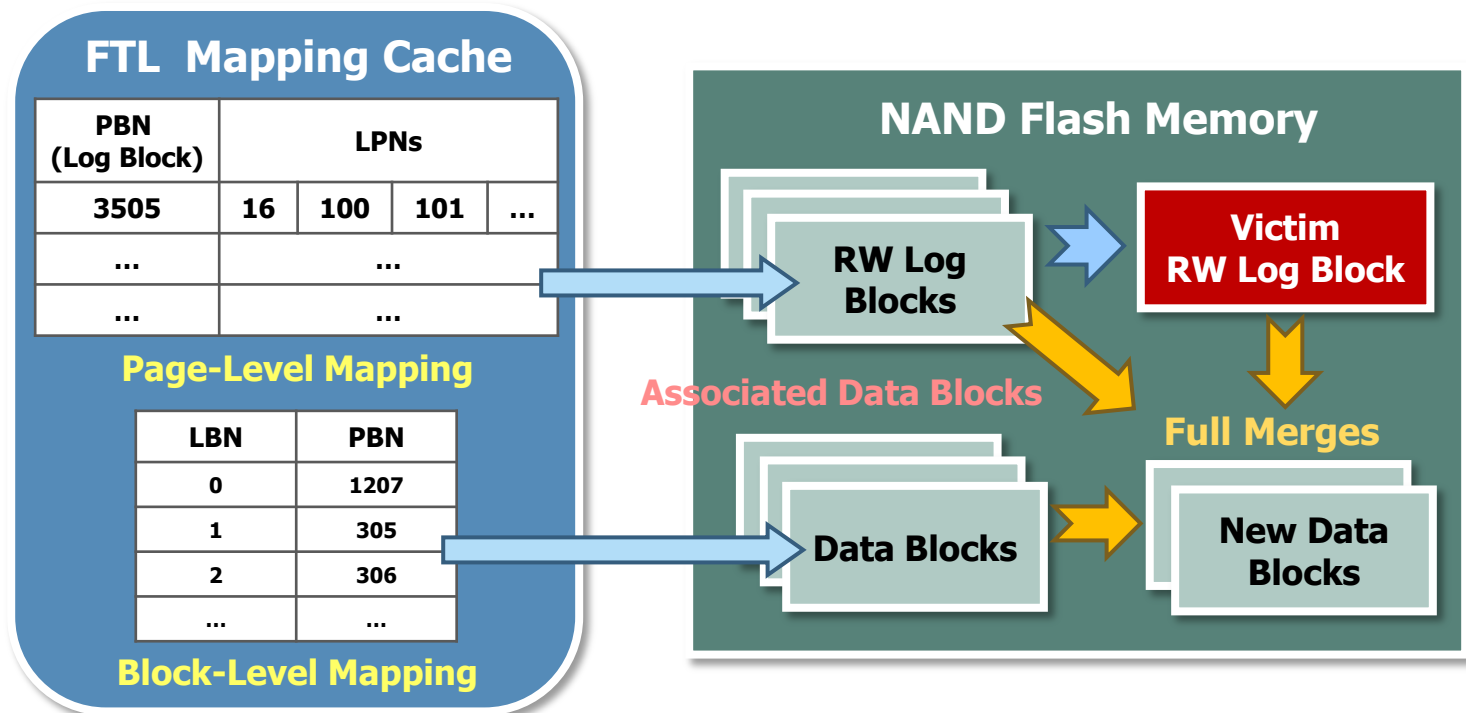


# Case Study II

## Fully Associative Sector Translation (FAST)

11/21

- **Data blocks** are managed by block-level mapping
  - Where all pages must be fully and sequentially written
- A fixed number of **log blocks** handle updates



# Case Study II

## Fully Associative Sector Translation (FAST)

12/21

- Write-dominant and high temporal-locality
  - Many valid pages in log blocks can be invalidated by following updates **with enough log blocks**

Page-Level Mapping

| PBN (Log Block) | LPNs |     |     |     |    |   |   |    |
|-----------------|------|-----|-----|-----|----|---|---|----|
| 3505            | 7    | 100 | 101 | 102 | 15 | 5 | 3 | 14 |
| ...             | ...  |     |     |     |    |   |   |    |

If all valid pages of the same associated data block (LBN:12) are invalidated, a full merge is avoided

- Read-dominant or small working set
  - Many log blocks remain unused, unnecessarily wasting the device cache
  - The unused mapping space can be utilized for buffering **by reducing # of available log blocks**

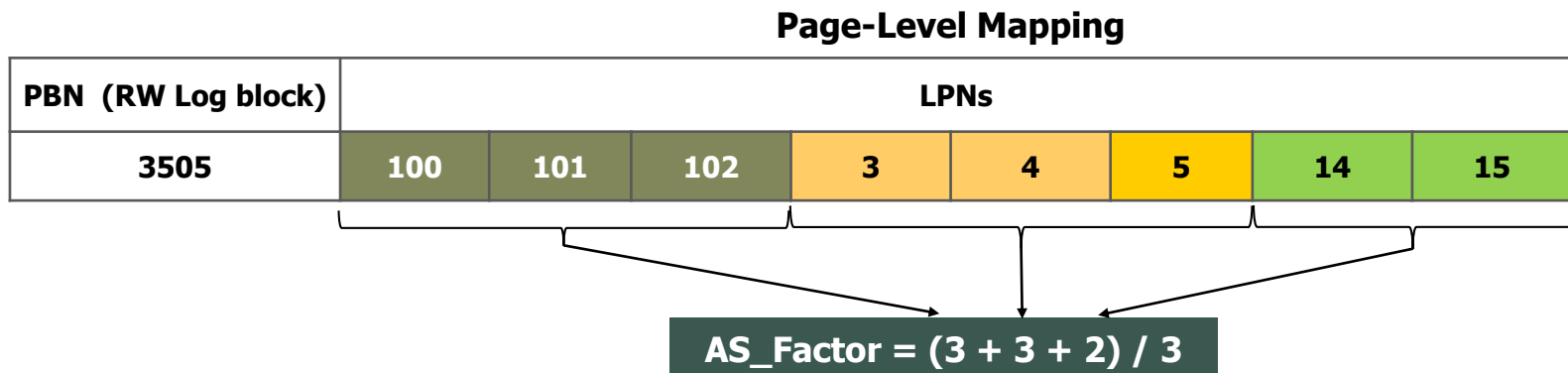
Page-Level Mapping

| PBN (Log Block) | LPNs  |     |     |     |
|-----------------|-------|-----|-----|-----|
| 3505            | 16    | 100 | 101 | ... |
| ...             | Empty |     |     |     |
| ...             | Empty |     |     |     |

# Opportunity Cost of the Mapping Cache with FAST

13/21

- Opportunity cost for a read miss in mapping cache
  - ▣ There is no read miss in FAST
- Opportunity cost for a write miss in mapping cache
  - ▣ *Full\_Merge\_Overhead / AS\_Factor*
  - ▣ *AS\_Factor* means the avg. # of written pages that belong to the same associated data block



# Opportunity Cost of the **Buffer Cache** with FAST

14/21

- **Opportunity cost for a read miss in buffer cache**
  - ▣ *Flash\_Read*
- **Opportunity cost for a write miss in buffer cache**
  - ▣ *Flash\_Write + RW\_Log\_Merge\_Cost / #Pages\_Per\_Block*

# Configurations for Experiments

15/21

- **8MB or 16MB of DRAM is assumed as the device cache**

|                         | <b>DFTL</b>           | <b>FAST</b>            |
|-------------------------|-----------------------|------------------------|
| <b>Tuning Interval</b>  | <b>1,000 Requests</b> | <b>10,000 Requests</b> |
| <b>Tuning Unit Size</b> | <b>A CMT Entry</b>    | <b>A Log Block</b>     |

- **64GB SLC NAND flash memory**
  - **The number of extra blocks is set as up to about 10% of the total capacity**

| <b>Flash Type</b> | <b>Unit Size (KB)</b> |              | <b>Access Time (<math>\mu</math>s)</b> |              |              |
|-------------------|-----------------------|--------------|--|--------------|--------------|
|                   | <b>Page</b>           | <b>Block</b> | <b>Read</b>                            | <b>Write</b> | <b>Erase</b> |
| <b>SLC</b>        | <b>2</b>              | <b>128</b>   | <b>72.8</b>                            | <b>252.8</b> | <b>1500</b>  |

# Summary of the Block-Level Traces

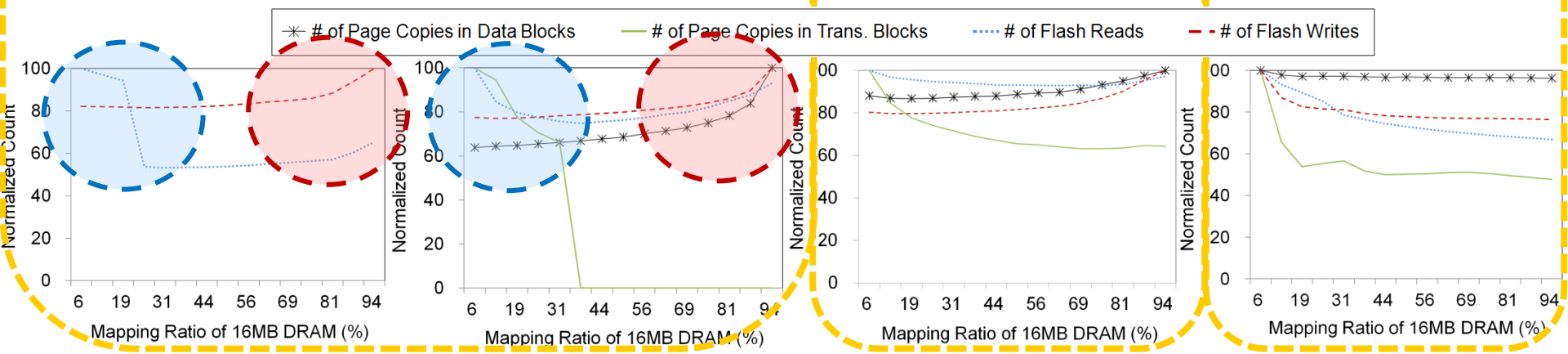
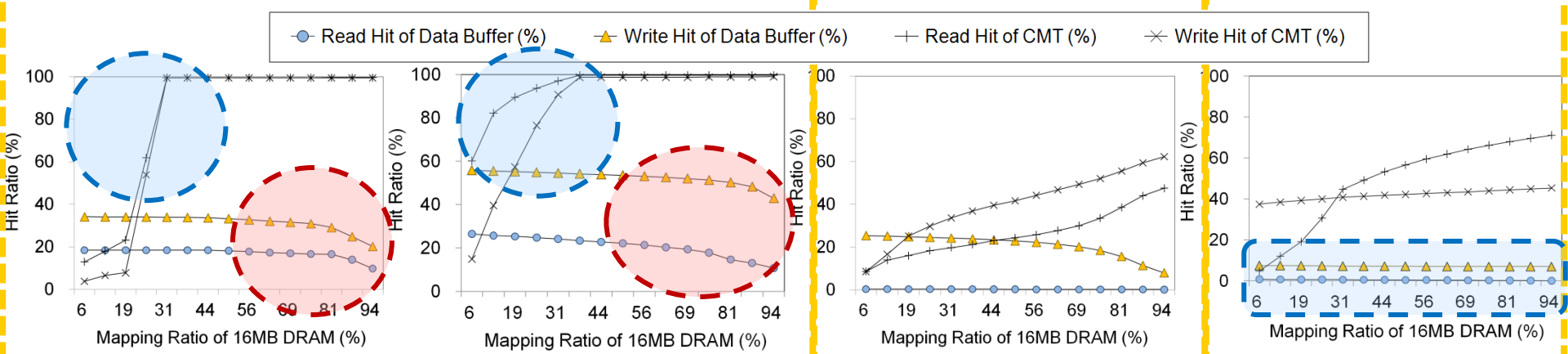
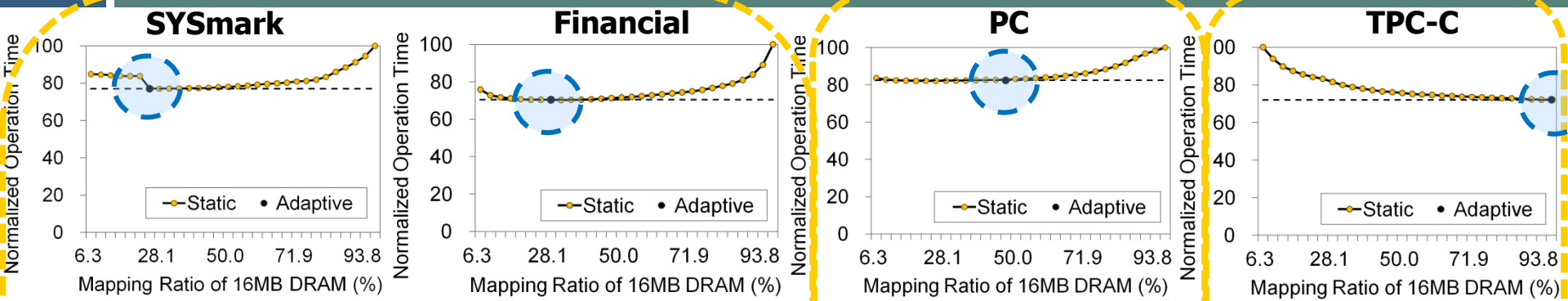
16/21

| Name             | Description   | Avg. Req. Size<br>[Read/Write]<br>(KB) | Req. Ratio<br>[Read/Write]<br>(%) | Working Set<br>[Read/Write]<br>(GB) |
|------------------|---|--|-----------------------------------|-------------------------------------|
| <b>SYSmark</b>   | Running <i>SYSmark 2007 Preview</i> including e-learning, office works, video creation, and 3D modeling | <b>13.6 / 20</b>                       | <b>33 / 67</b>                    | <b>0.11 / 0.24</b>                  |
| <b>Financial</b> | I/O trace from an OnLine Transaction Processing (OLTP) application running at a financial institution   | <b>2.3 / 3.6</b>                       | <b>47.4 / 52.6</b>                | <b>0.45 / 0.5</b>                   |
| <b>PC</b>        | Document-based realistic workloads using various office applications                                    | <b>20 / 13.4</b>                       | <b>23.7 / 76.3</b>                | <b>5.82 / 8.45</b>                  |
| <b>TPC-C</b>     | Running a TPC-C benchmark test with <i>Benchmark Factory</i>  | <b>2.2 / 2.1</b>                       | <b>81.4 / 18.6</b>                | <b>8.04 / 4.45</b>                  |



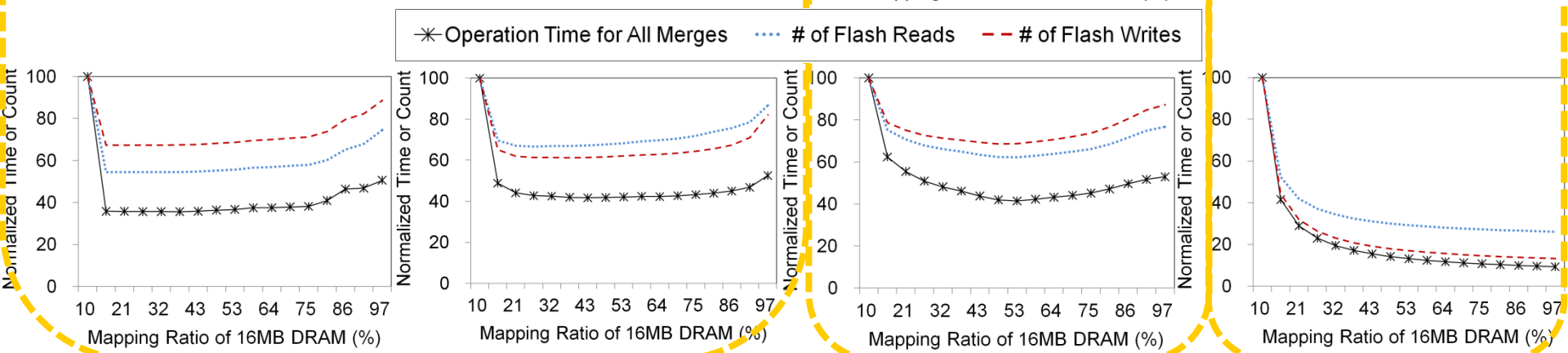
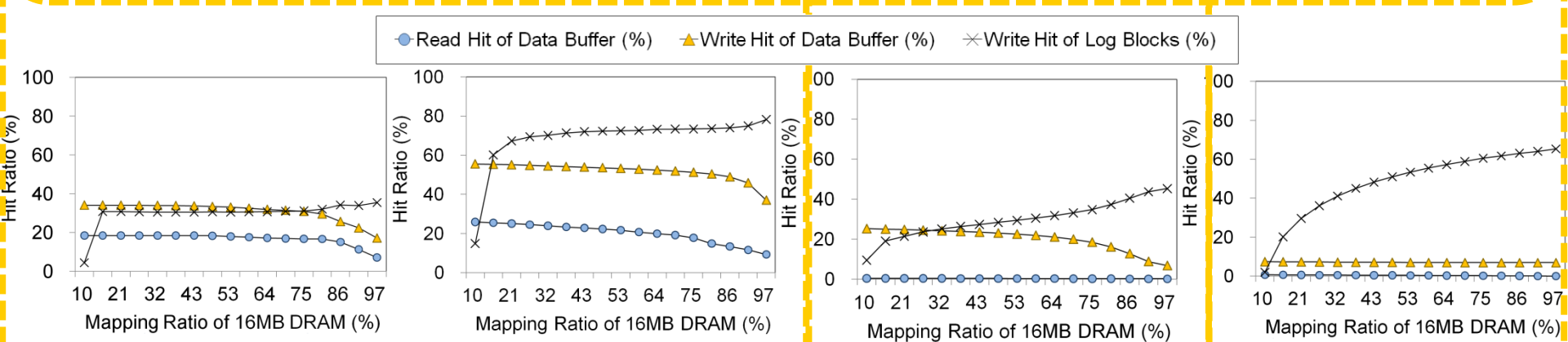
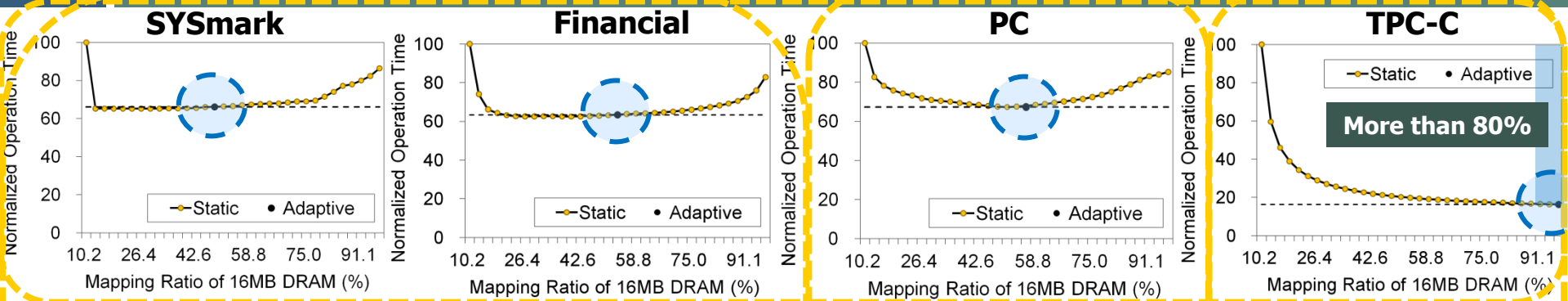
# Operation Time with DFTL

17/21



# Operation Time with FAST

18/21

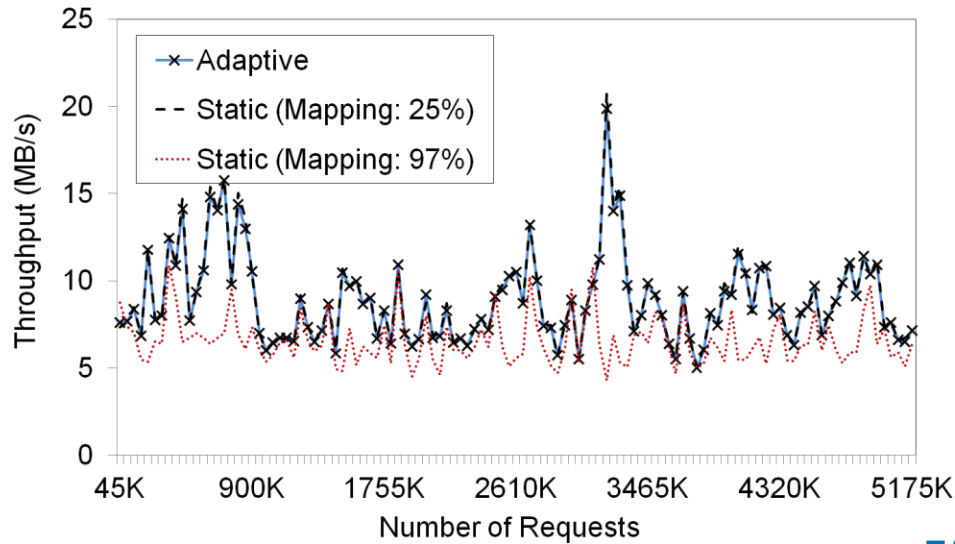


# Throughput with DFTL and FAST

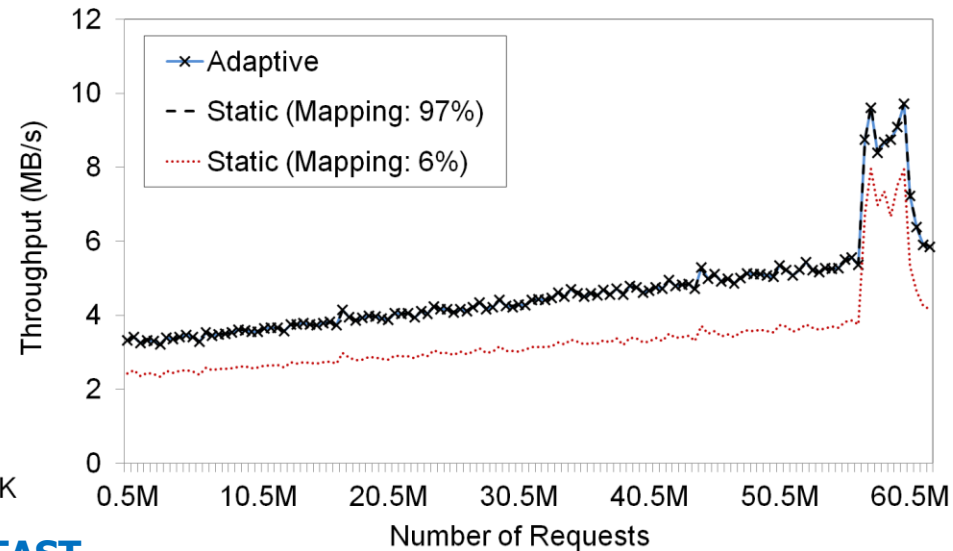
19/21

## DFTL

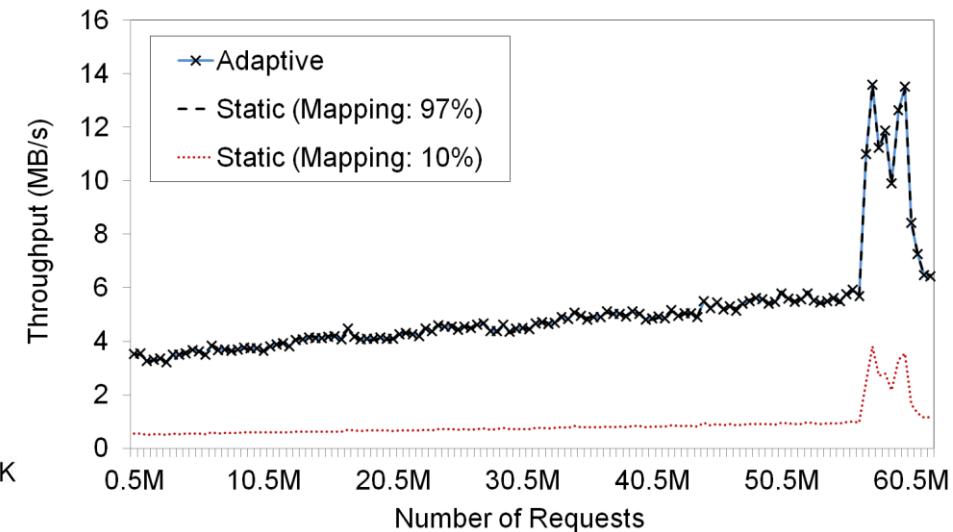
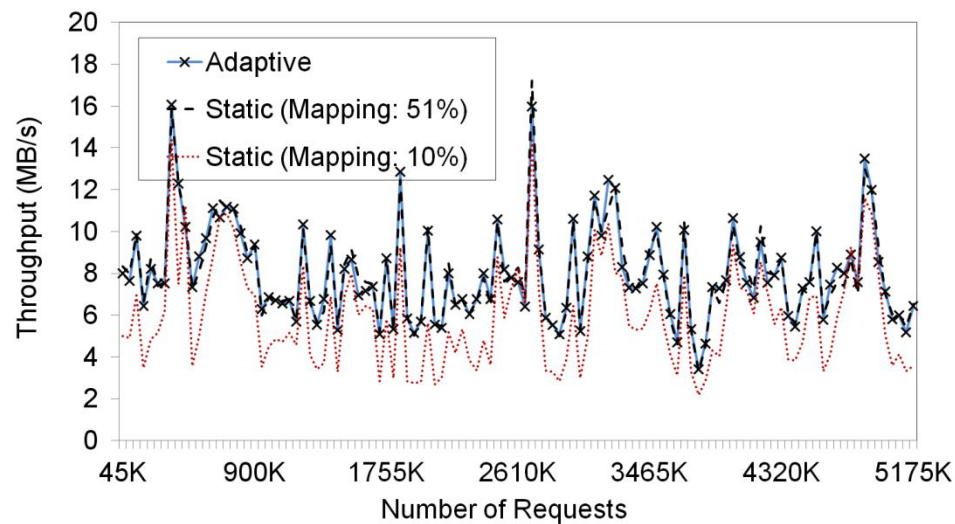
### PC



### TPC-C



## FAST

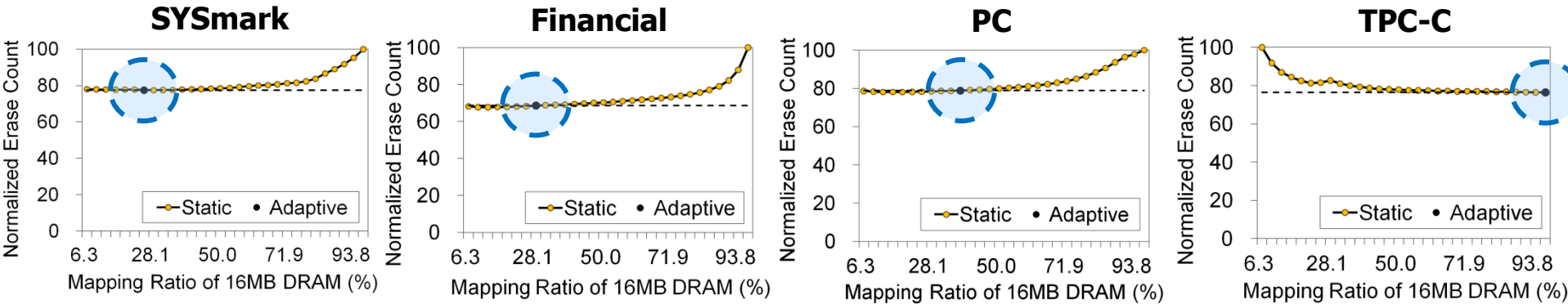


# Erase Count with DFTL and FAST

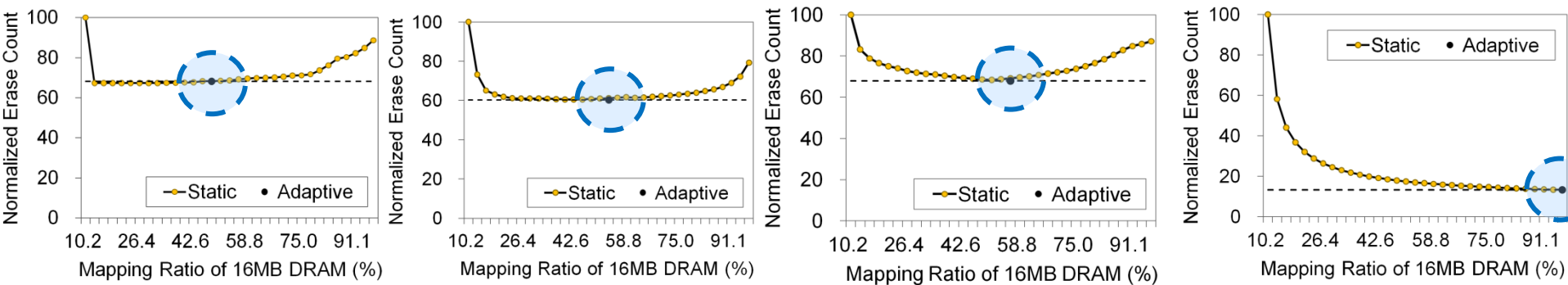
20/21

- The proposed scheme can extend the lifetime of SSDs

## DFTL



## FAST



# Conclusions

21/21

- We proposed an **adaptive partitioning scheme** for better performance of SSDs
  - The proposed scheme adaptively tunes the BM ratio according to workload characteristics
  - We built a **cost-benefit model** based on a *ghost caching mechanism*
  - The performance results come near the best performance under the static partitioning policy with varied workloads
- We expect that SSDs equipped with the proposed scheme can be deployed in different environments without workload-specific tuning

# Extra Slides



# Implementation of Ghost Mapping Cache

23/21

## □ *Bloom Filter*

- To insert an LPN, the corresponding hash bucket is set
  - If a hash collision occurs, the bit count of hash bitmap does not increase
- To flush a victim, we reset of the bit in random order
  - the bitmap hash is flushed only if **(the current bit count of the hash bitmap) \*  $\alpha$**  is smaller than the bit count that should be preserved without collisions
  - $\alpha$  was set as 3 for all the simulation results